

Storage Evaluation on FG, FC, and GPCF

Overview

- Introduction
- Lustre Evaluation: Status
 - The test bed
 - Results from std benchmarks
 - Results from app-based benchmarks
- Metrics measurement for the RunII experiments
- Conclusions and future work

Jun 29, 2010

Gabriele Garzoglio
Computing Division, Fermilab

Context

- Goal
 - Evaluation of storage technologies for the use case of data intensive Grid jobs.
- Technologies considered
 - Hadoop Distributed File System (HDFS)
 - Lustre
 - Blue Arc (BA)
- Targeted infrastructures:
 - FermiGrid, FermiCloud, and the General Physics Computing Farm.

Collaboration

- REX was asked by Patty to do an evaluation of storage technologies looking into the future
- REX and DOCS have agreed to work together to reuse the infrastructure at FermiCloud allocated to the original evaluation project AND take advantage of REX contacts with the users
- Other partners in the evaluation include the FermiGrid / FermiCloud, OSG Storage, DMS, and FEF groups at Fermilab

Evaluation Method

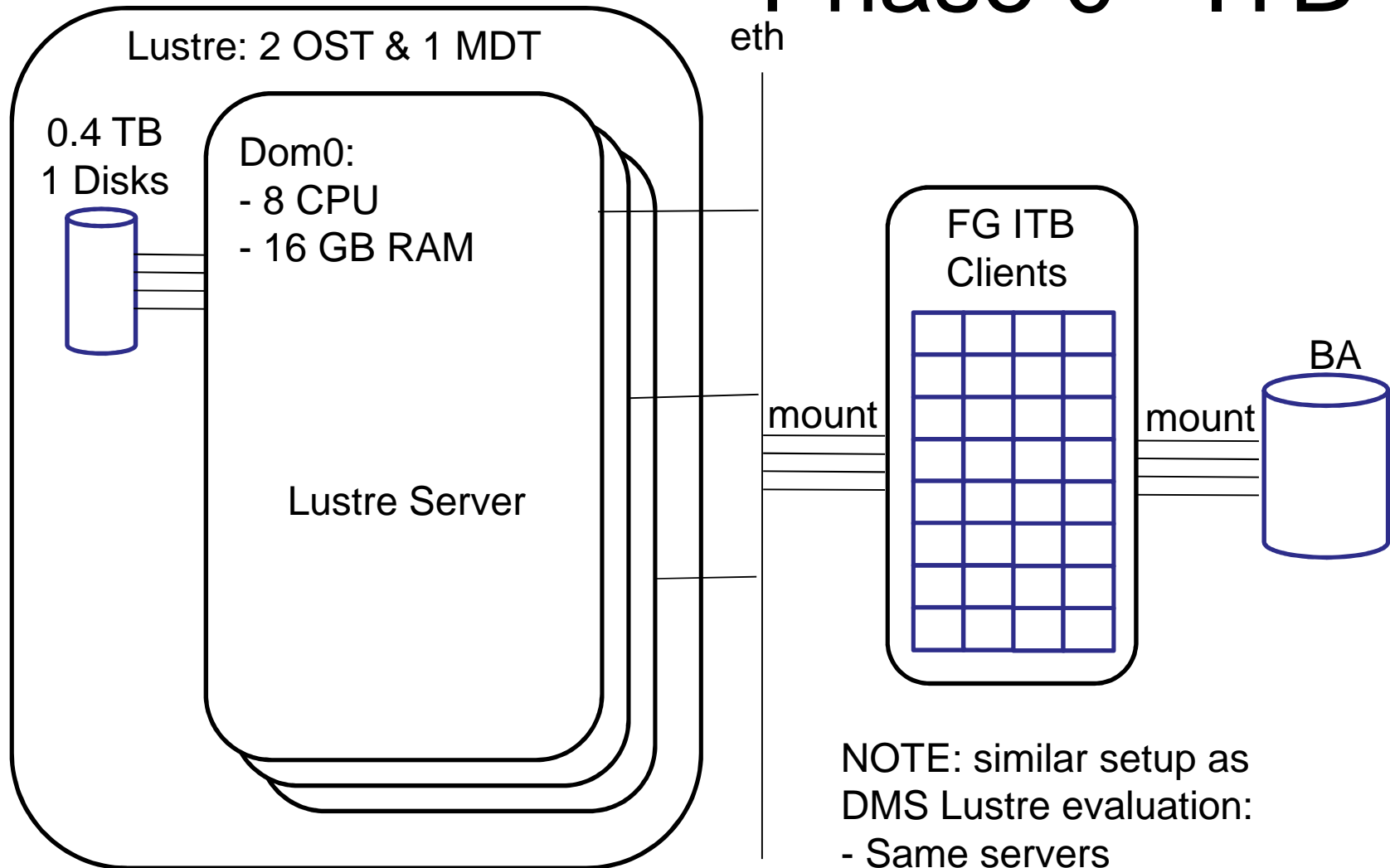
- **Set the scale:** measure storage metrics from running experiments to set the scale on expected bandwidth, typical file size, number of clients, etc. (Done)
- **Check sanity:** run standard benchmarks on storage installations
- **Measure performance:** study response of the technology to real-life applications access patterns
- **Fault tolerance:** simulate faults and study reactions

Machine Specifications

- FermiCloud (FC) Server Machines:
 - Lustre 1.8.3: Striped across 3 OSS, 1 MB block
 - CPU: dual, quad core Xeon E5640 @ 2.67GHz with 12 MB cache, 24 GB RAM
 - Disk: 6 SATA disks in RAID 5 for 2 TB + 2 sys disks
(hdparm on raid → 376.94 MB/sec)
 - 1 GB Eth + IB (not used yet)
- ITB Server Machines:
 - Lustre 1.8.3 : Striped across 2 OSS, 1 MB block
 - CPU: dual, quad core Xeon X5355 @ 2.66GHz with 4 MB cache: 16 GB RAM
 - Disk: single 500 GB disk
(hdparm on disk → 76.42 MB/sec)
- DMS Server Machine: Lustre 1.6
 - same specs as ITB
- ITB Client Machines
 - same specs as ITB

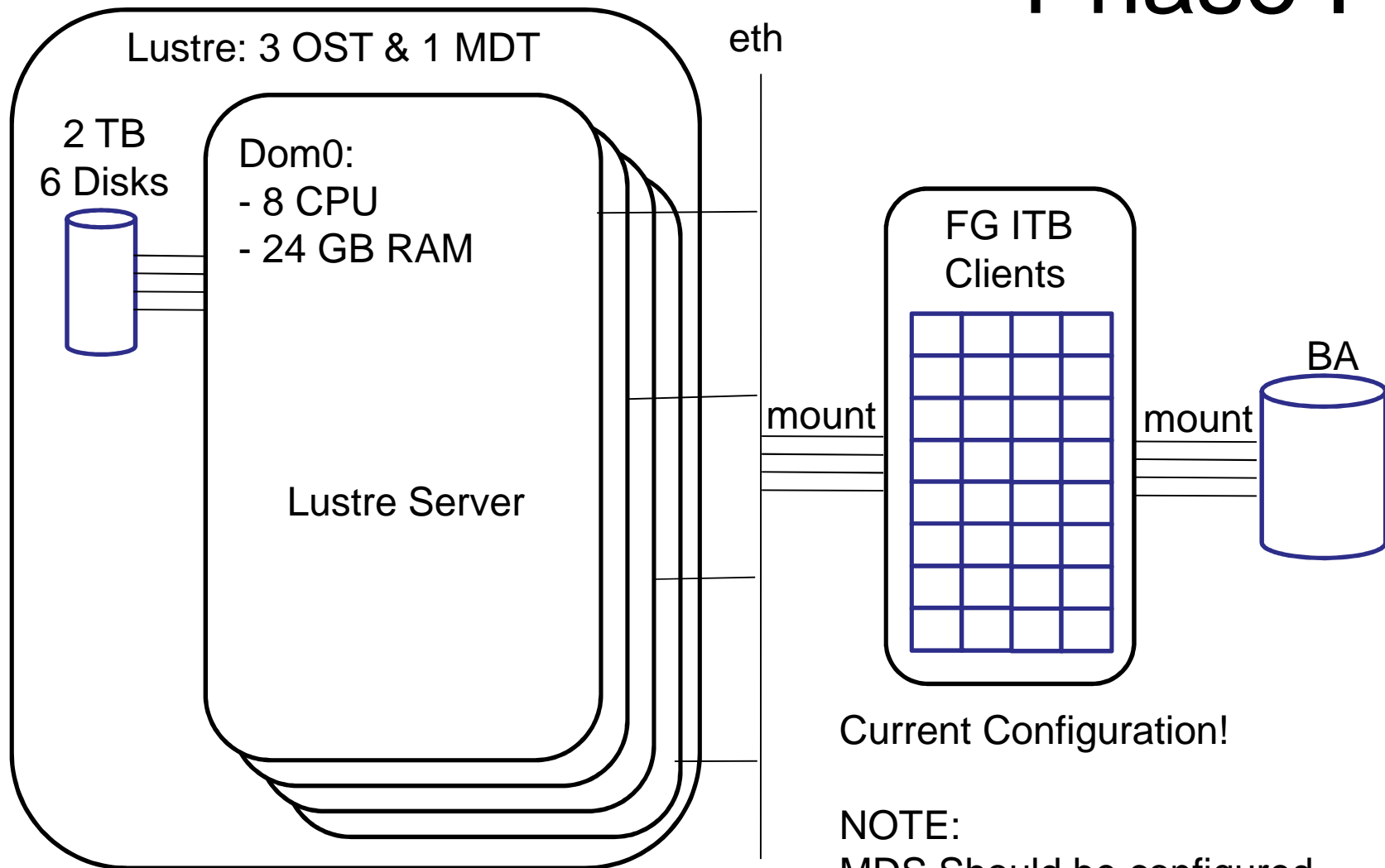
Lustre Evaluation Test Bed

Phase 0 - ITB



NOTE: similar setup as
DMS Lustre evaluation:
- Same servers
- 2 OST vs. 3 OSG for DMS.

Lustre Evaluation Test Bed Phase I

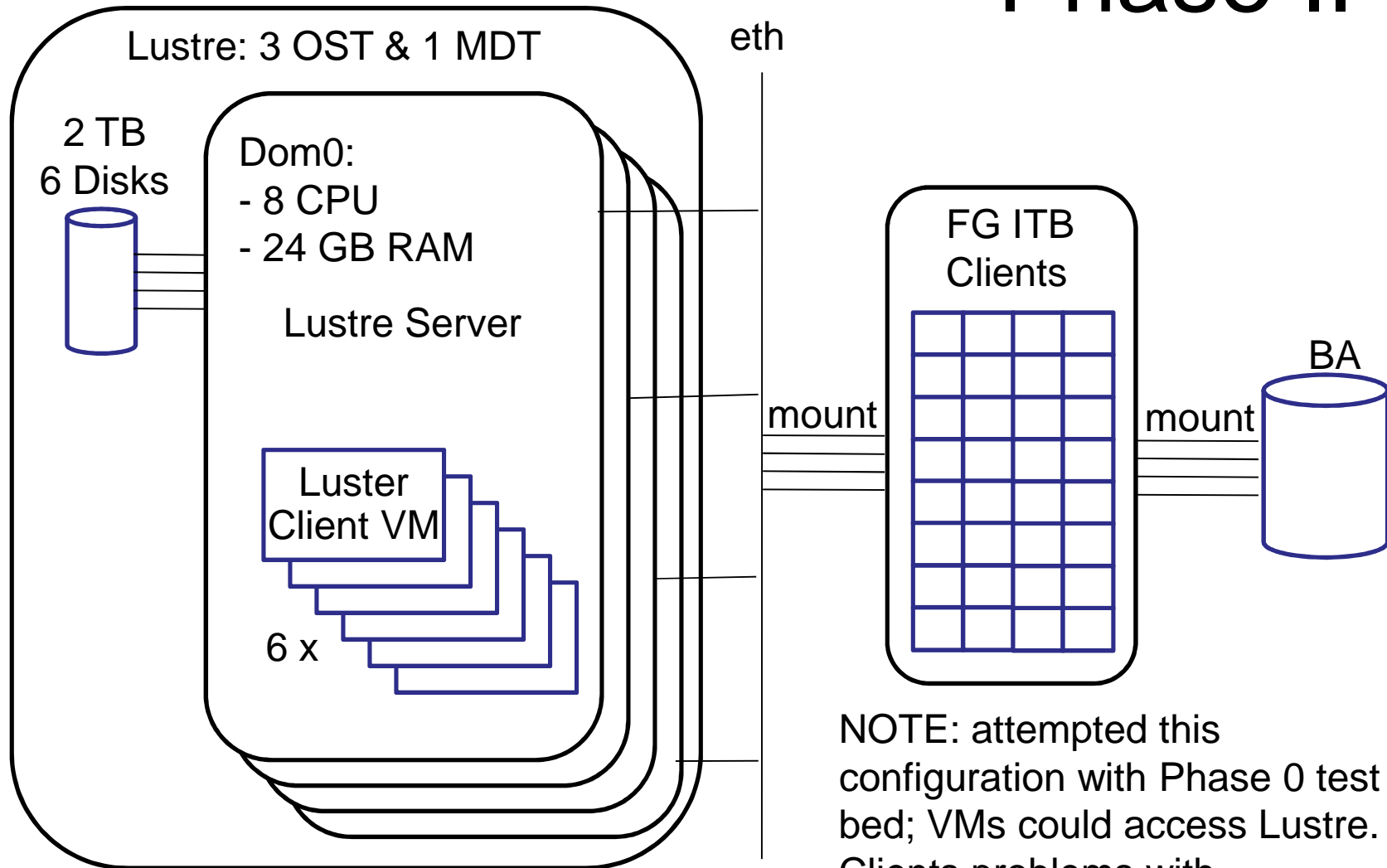


Current Configuration!

NOTE:

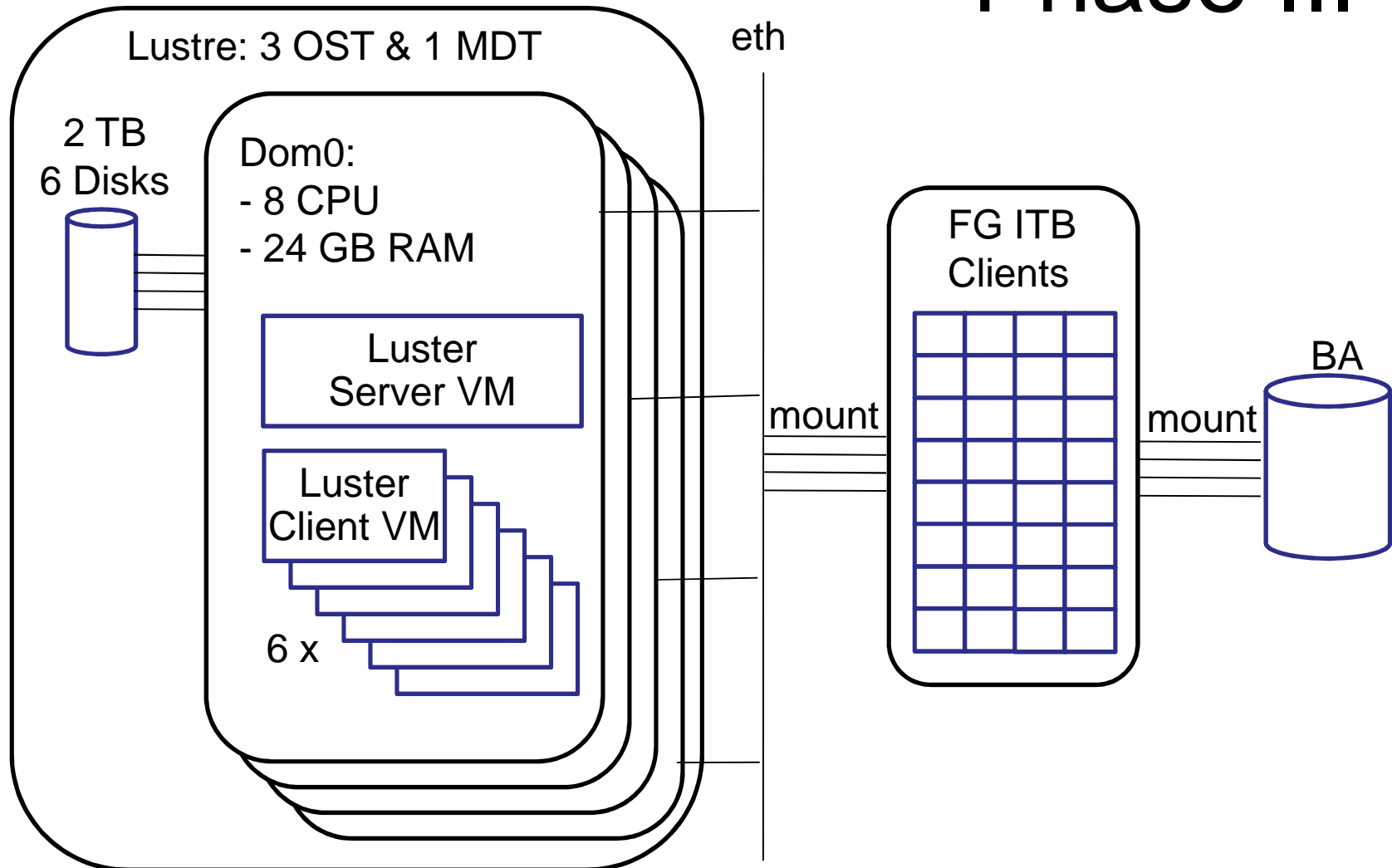
MDS Should be configured
as RAID10 rather than RAID 5

Lustre Evaluation Test Bed Phase II



NOTE: attempted this configuration with Phase 0 test bed; VMs could access Lustre. Clients problems with paravirtualizer.

Lustre Evaluation Test Bed Phase III



Standard Benchmarks

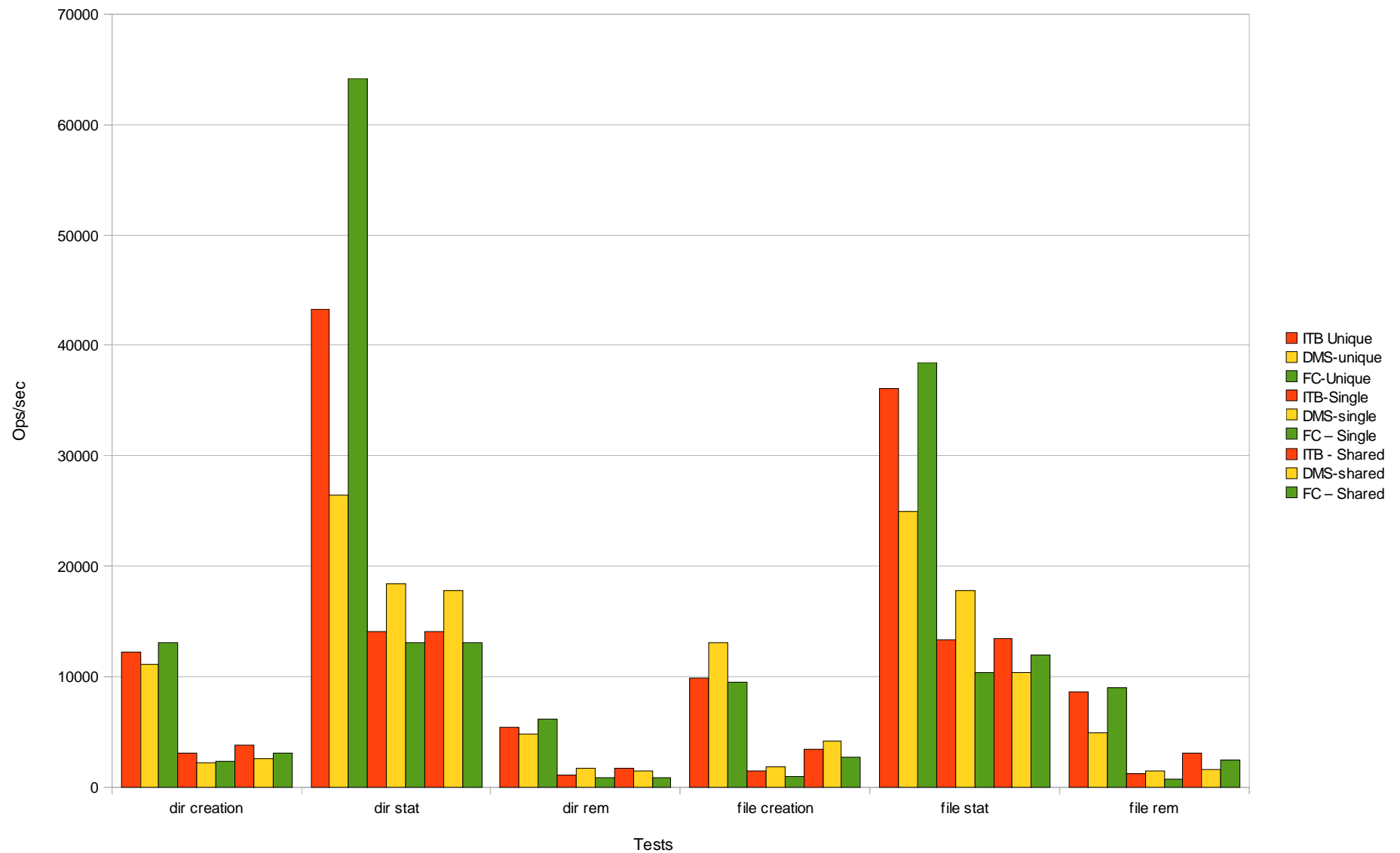
Results - Summary

Mdtest – Tests metadata rates from multiple clients. File/Directory Creation, Stat, Deletion. **Setup:** 48 clients on 6 VM/nodes. **Result:** Within 25% or faster of DMS results.

IOZone – Writes (2GB) file from each client and performs read/write tests. **Setup:** 3-48 clients on 3 VM/nodes. **Result:** Write results match DMS report, ~70MB/sec ITB, ~85MB/sec DMS, ~110MB/sec FC. Read results vary based on memory/file-size due to caching effects

Fileop – IOZone's metadata tests. Tests rates of mkdir, chdir, open, close, etc. **Setup:** 1 client. **Result:** Within 30% of DMS results on most tests.

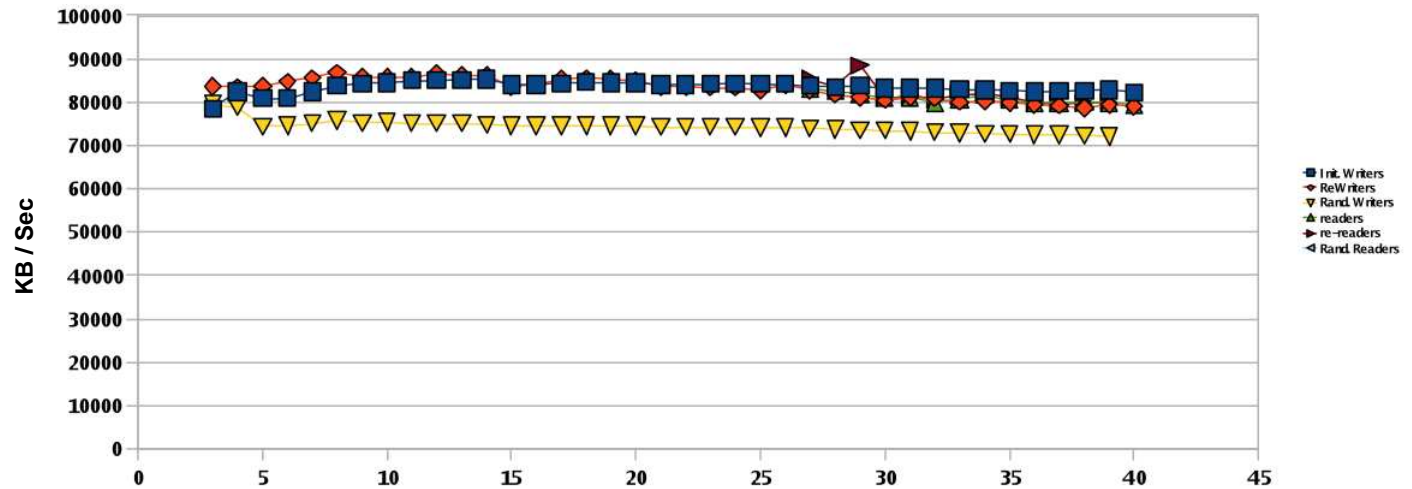
Metadata Test Results



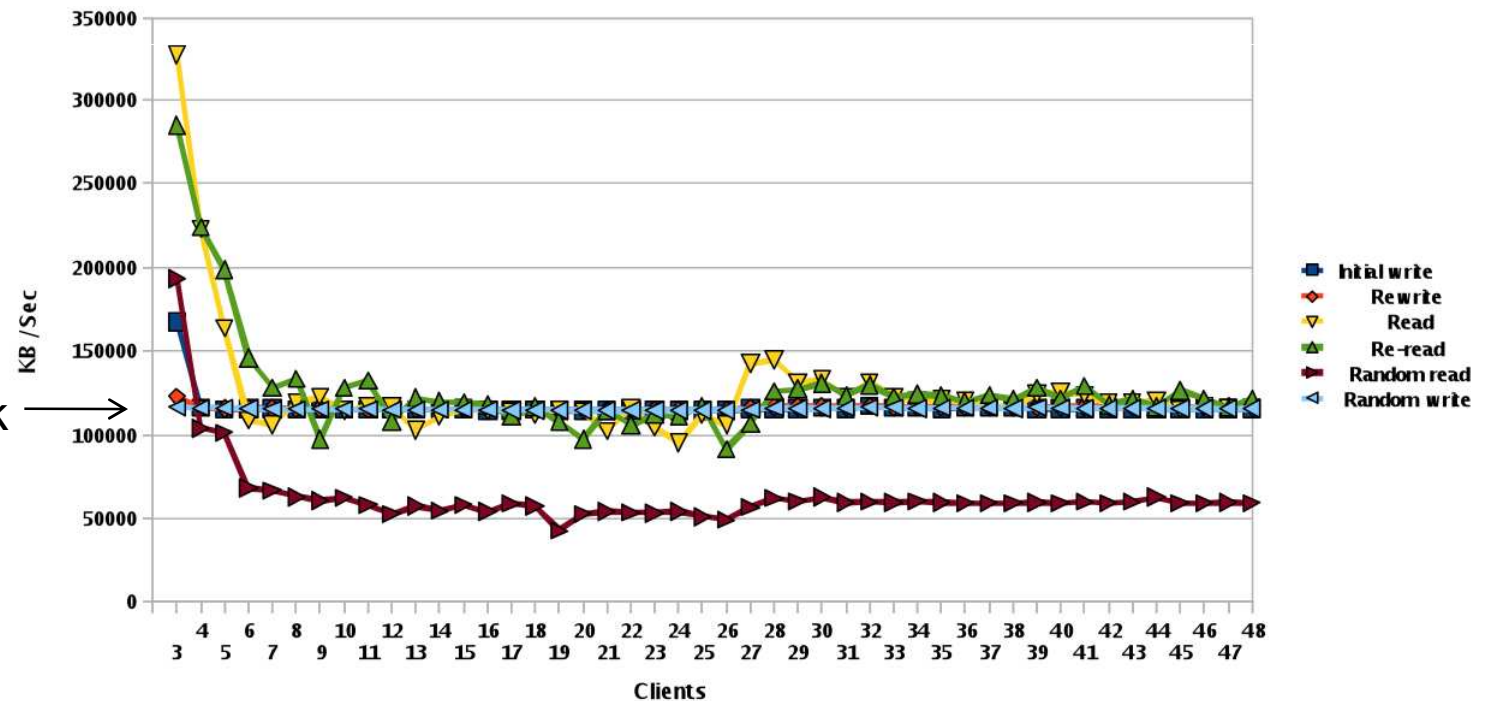
48 clients on 6 VM on 6 different nodes

Storage Evaluation on FG, FC, and GPCF

DMS Evaluation



Fermicloud IO Zone Results



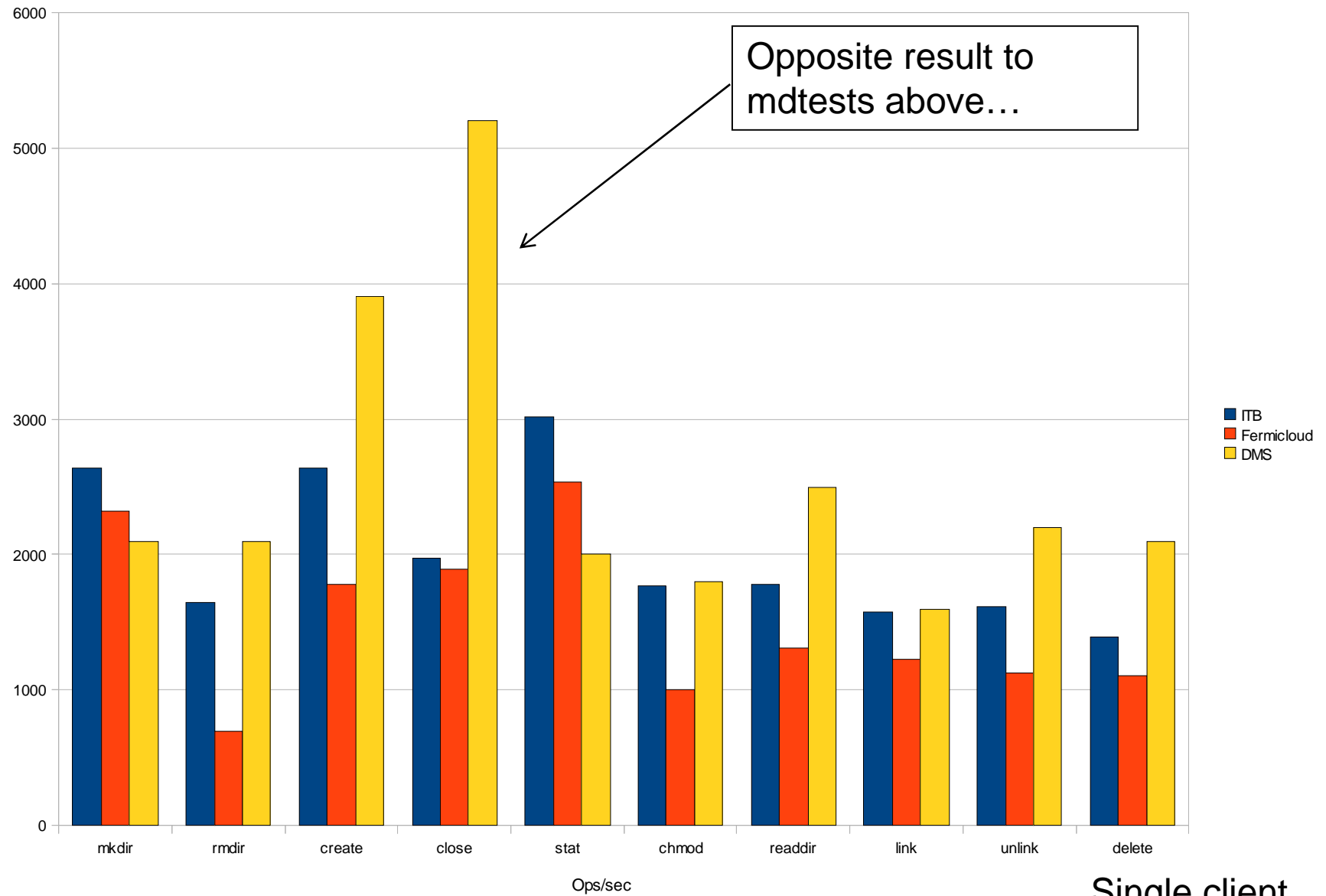
Up to 48 clients
on 3 VM on 3
different nodes

Plots show
aggregate
bandwidth

1 Gbit / s
Possible Network
Bottleneck...

Storage Evaluation on FG, FC, and GPCF

Fileop Results

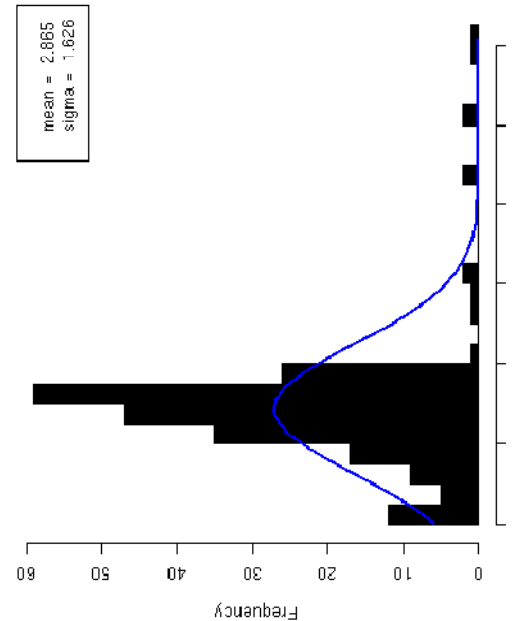


Response to real-life apps

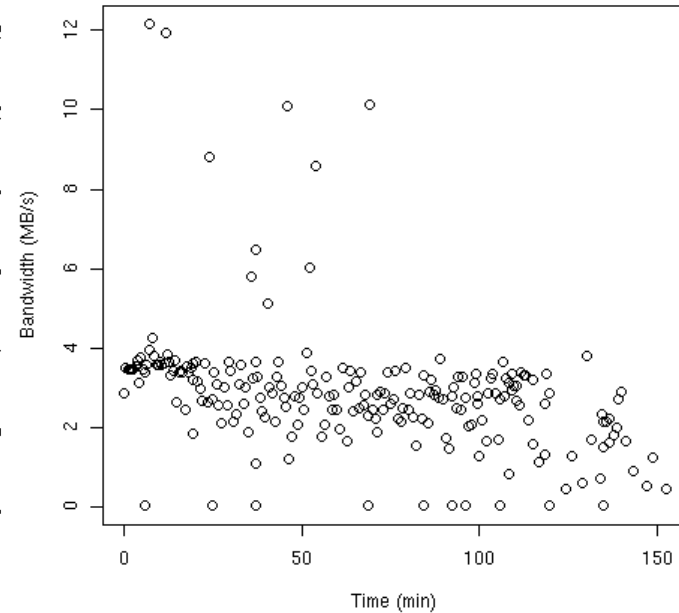
- Lustre is pre-loaded with root data from Minos and Nova (mc)
- A job invokes the application to run on 1 file for all files sequentially
- Multiple clients are run simultaneously (1 or 9 or 21)
- We measure bandwidth to storage, file size, and file processing time.

Minos

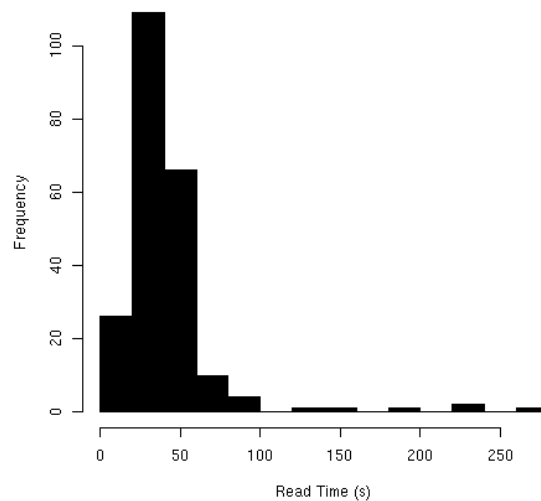
- 21 Clients
- Minos application (loon) skimming
- Random access to 1400 files (interrupted)
- Access is CPU-bound (checked via top)



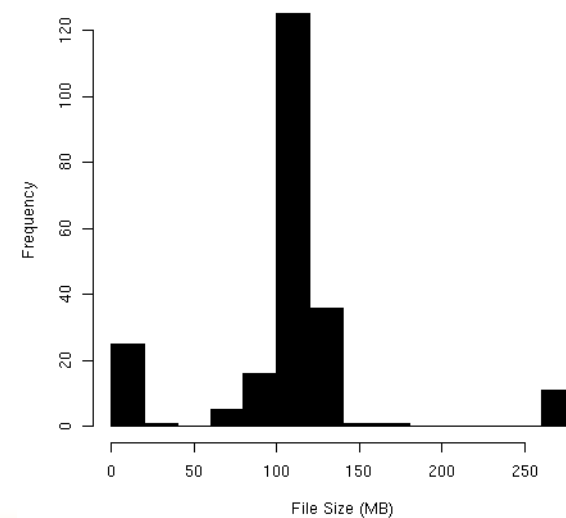
Bandwidth with 21 minos clients - Rand access
FC Lustre



Read time distribution - Rand access - 21 minos clients
FC Lustre

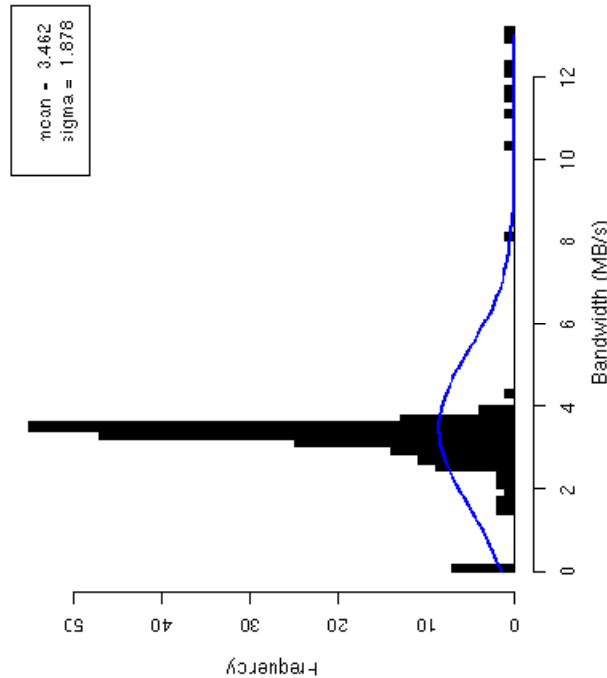


File Size distribution - Rand access - 21 minos clients
FC Lustre

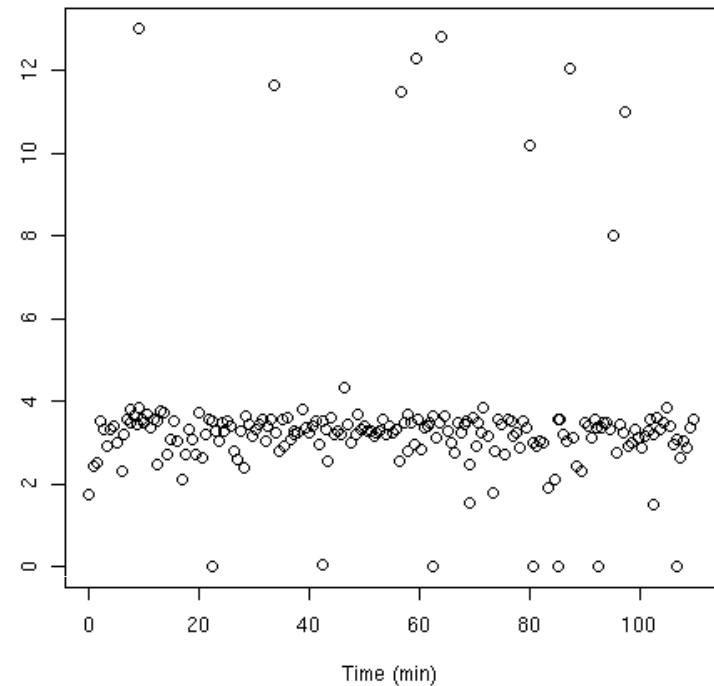


Minos

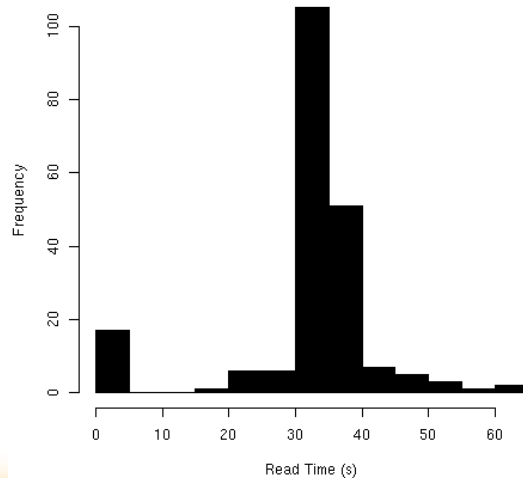
- 9 Clients
- Minos application (loon) skimming
- Random access to 1400 files (interrupted)
- Access is CPU-bound



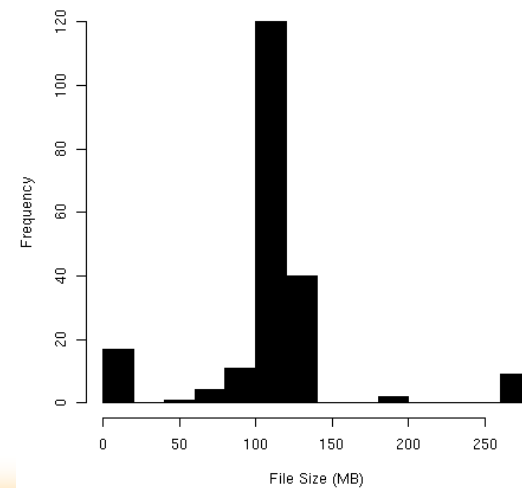
Bandwidth with 9 minos clients - Rand access
FC Lustre



Read time distribution - Rand access - 9 minos clients
FC Lustre

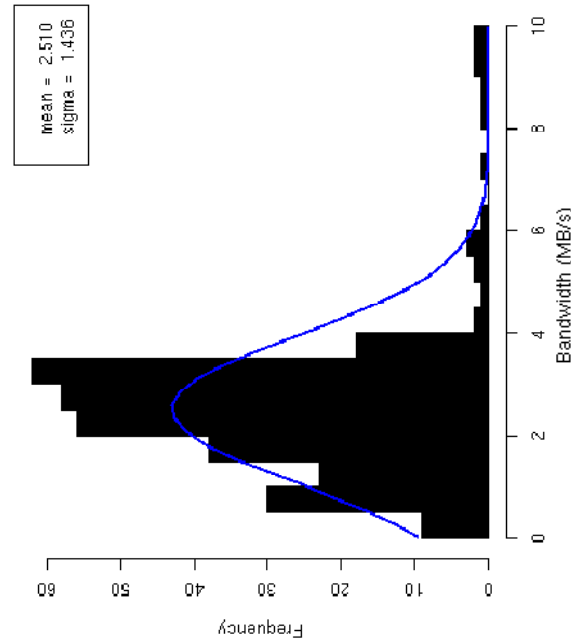


File Size distribution - Rand access - 9 minos clients
FC Lustre

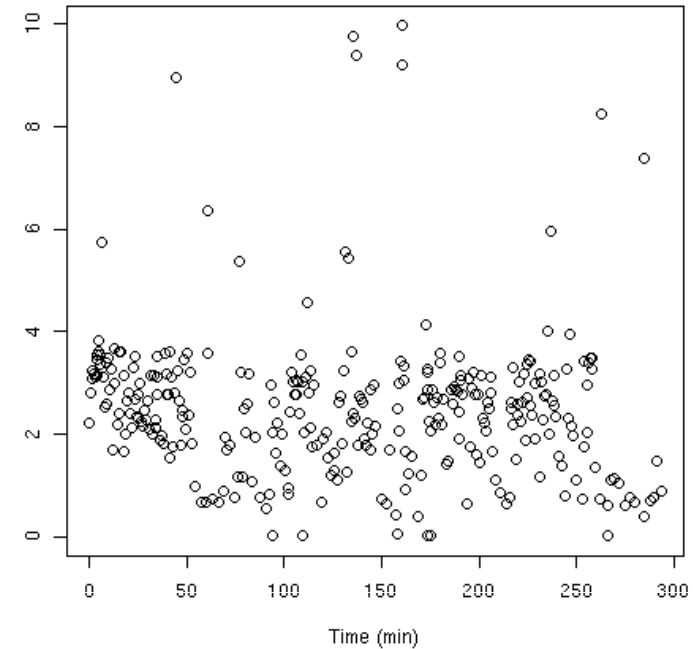


Minos

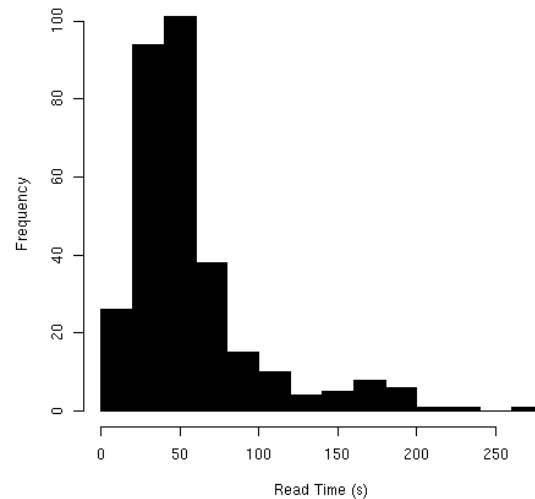
- 21 Clients
- Minos application (loosely) skimming
- Seq. access to 1400 files (interrupted)
- Access is CPU-bound



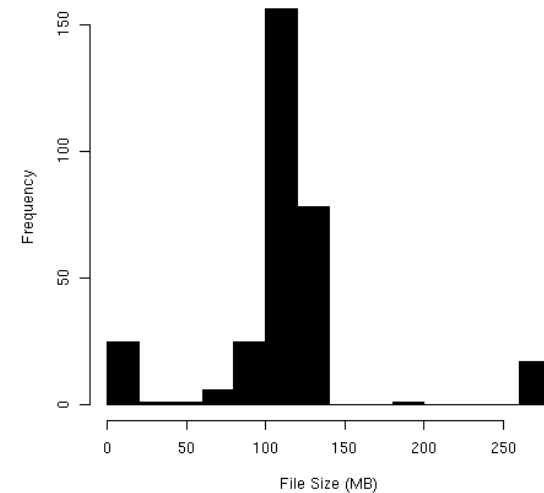
Bandwidth with 21 minos clients - Seq access
FC Lustre



Read time distribution - Seq access - 21 minos clients
FC Lustre

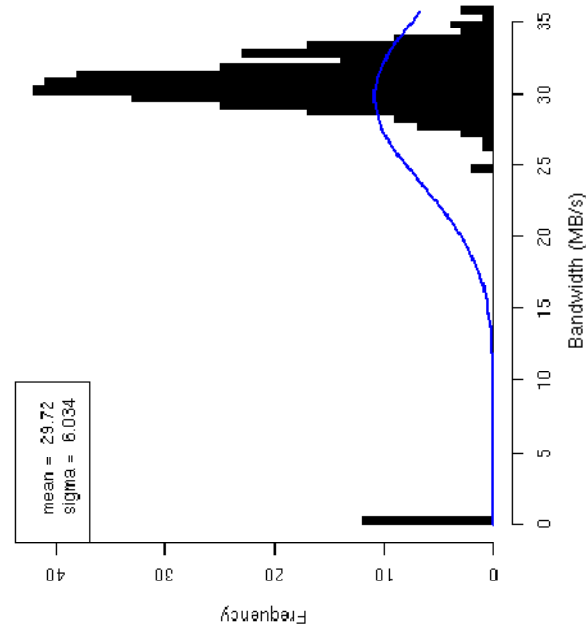


File Size distribution - Seq access - 21 minos clients
FC Lustre

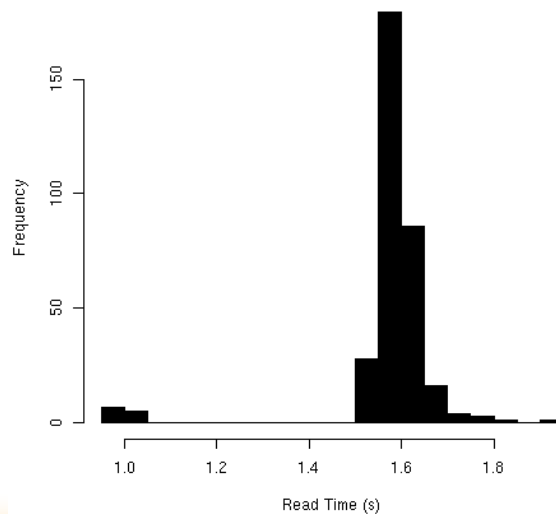


Nova

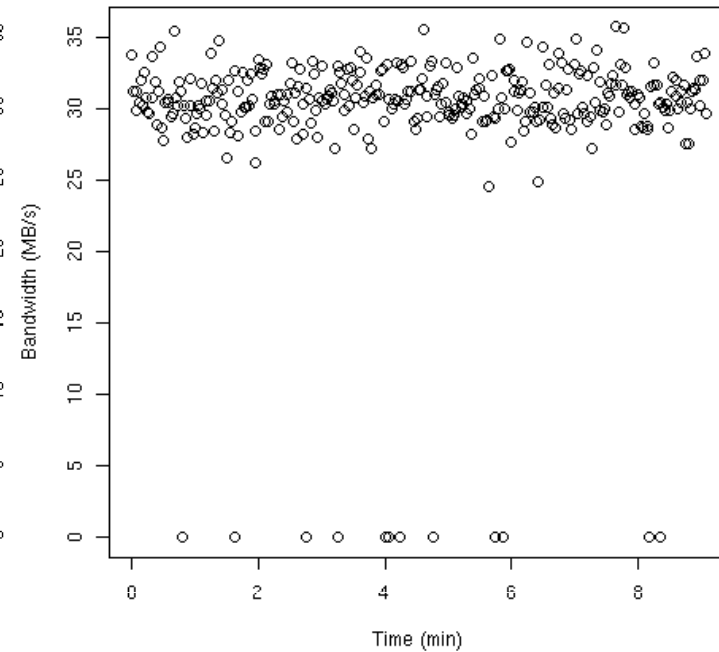
- 21 Clients
- Nova application (ana) skimming
- Rand access to 330 files:
 real 9m5.8s
 user 5m11.4s
 sys 1m28.4s



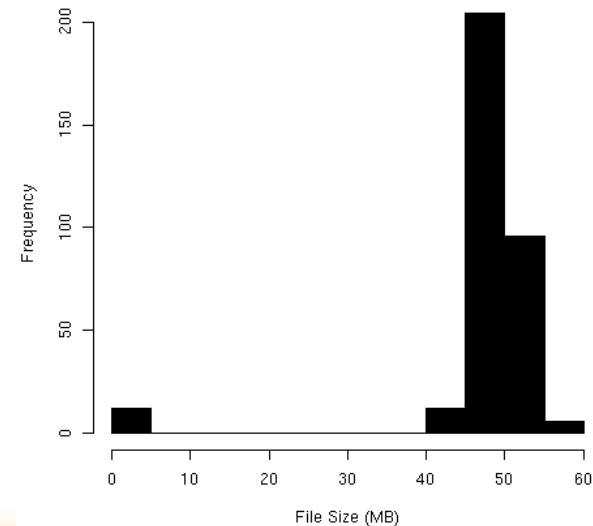
Read time distribution - Rand access - 21 nova clients
FC Lustre



Bandwidth with 21 nova clients - Rand access
FC Lustre



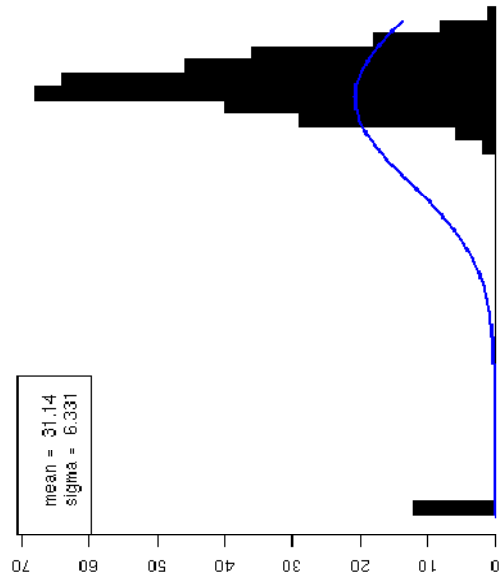
File Size distribution - Rand access - 21 nova clients
FC Lustre



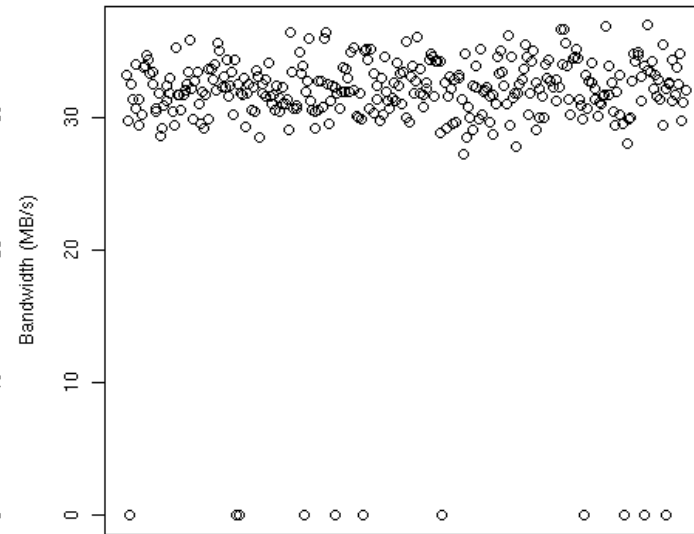
Nova

- 9 Clients
- Nova application (ana) skimming
- Rand access to 330 files:

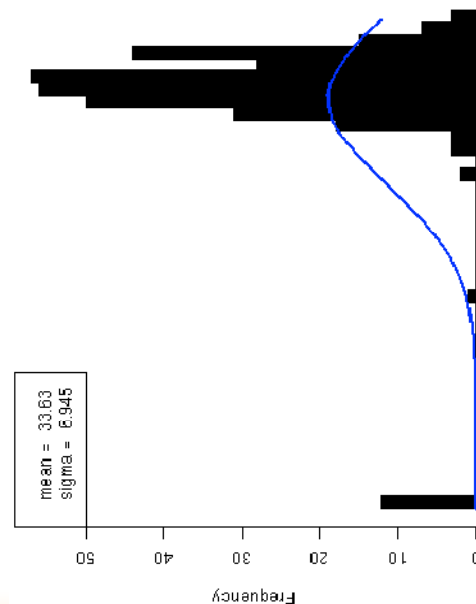
real 8m54.1s
user 5m11.9s
sys 1m28.2s



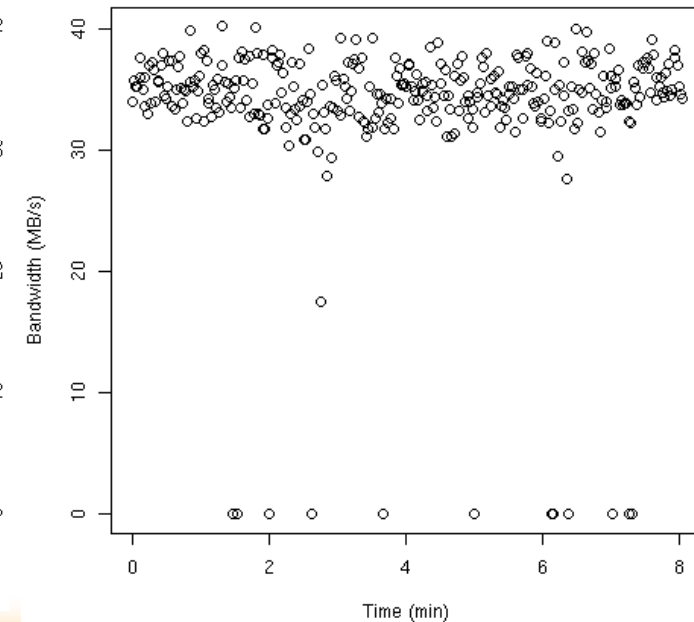
Bandwidth with 9 nova clients - Rand access
FC Lustre



- 1 Client
- Nova application (ana) skimming
- Rand access to 330 files

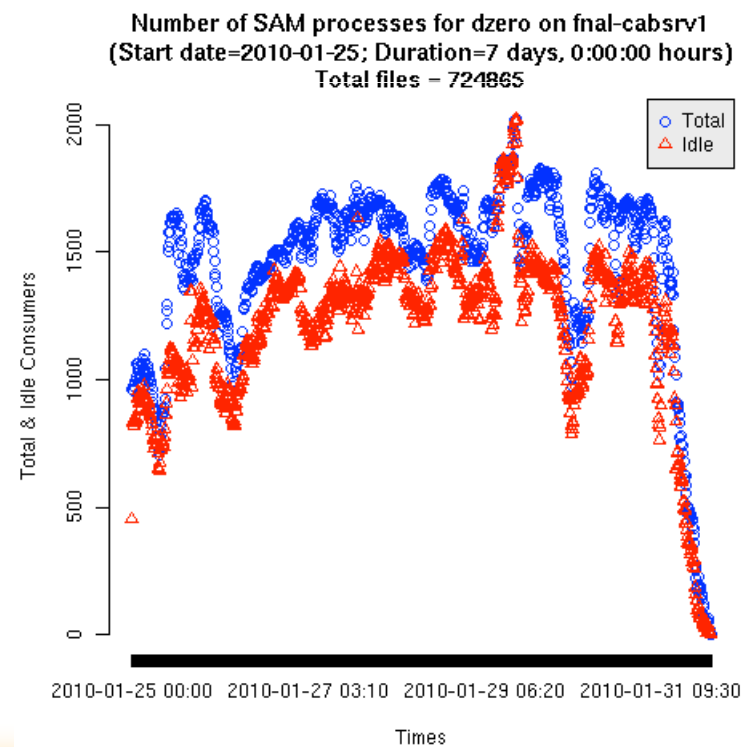
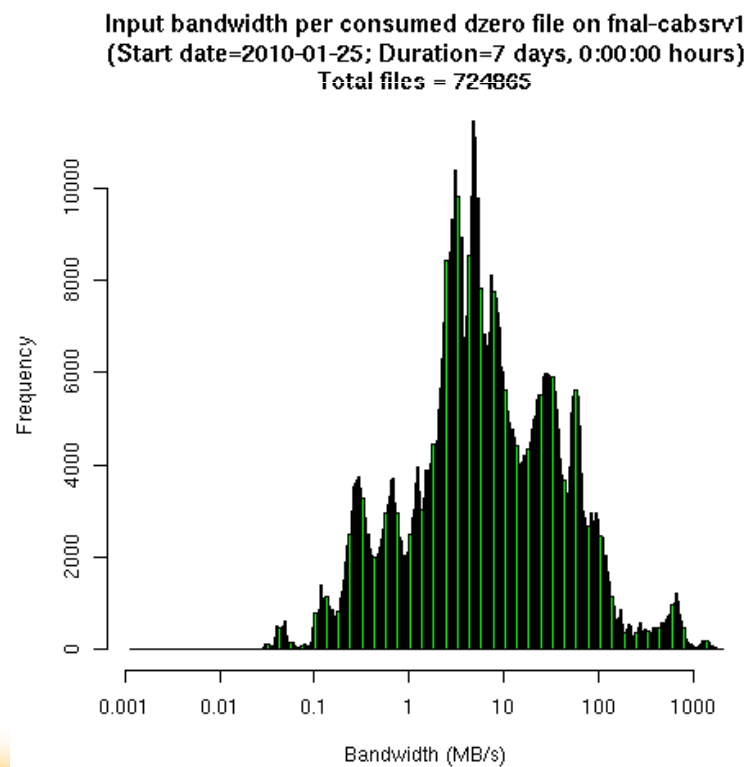
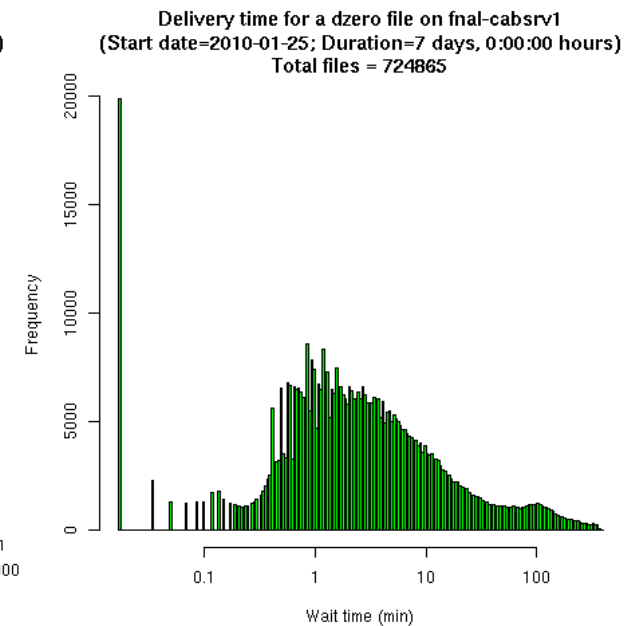
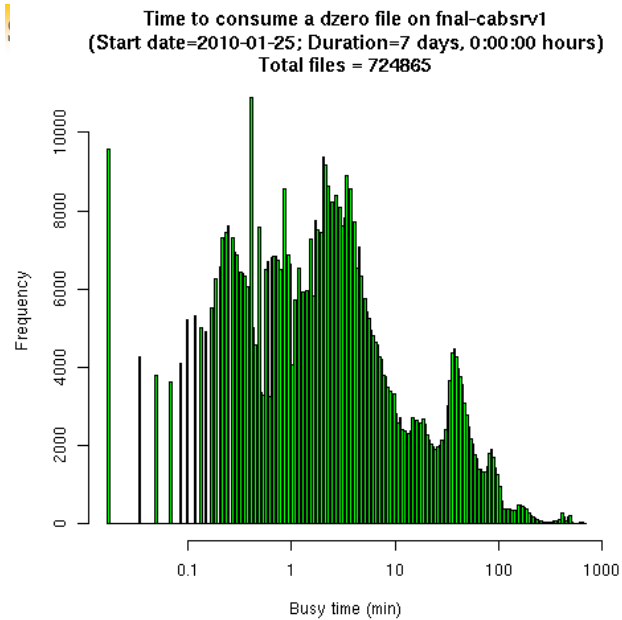
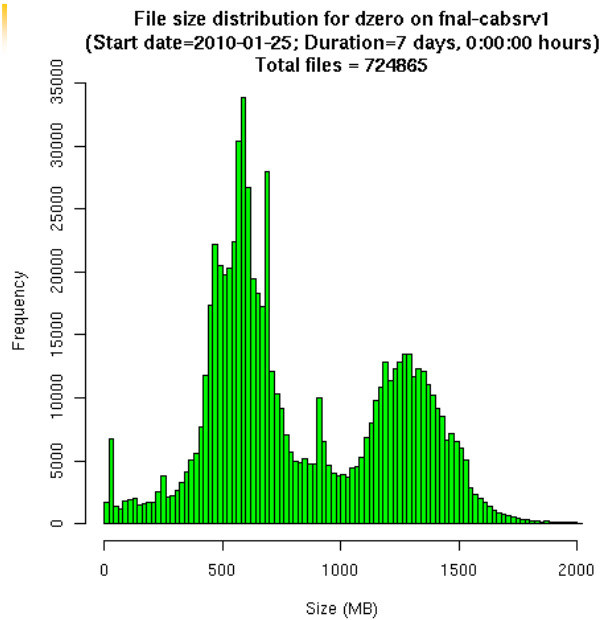


Bandwidth with 1 nova client - Rand access
FC Lustre



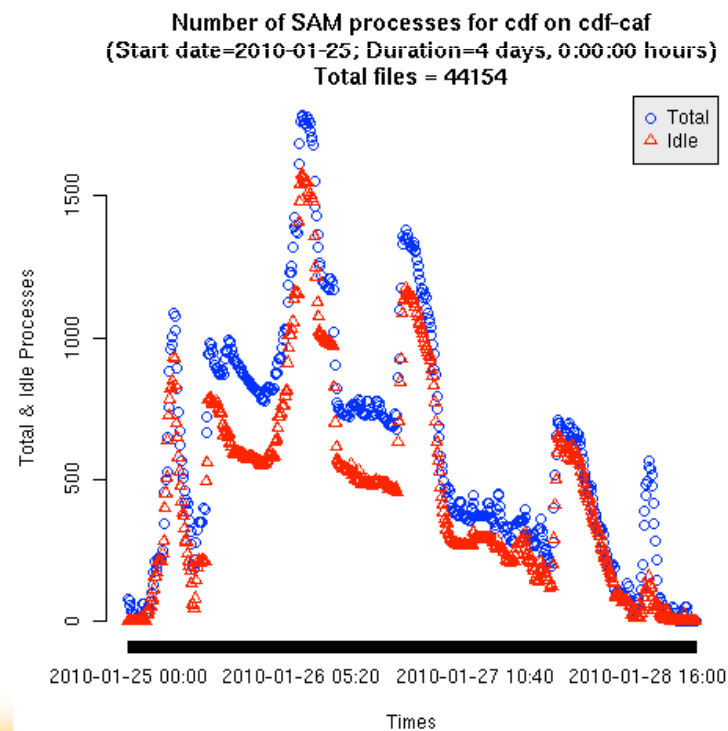
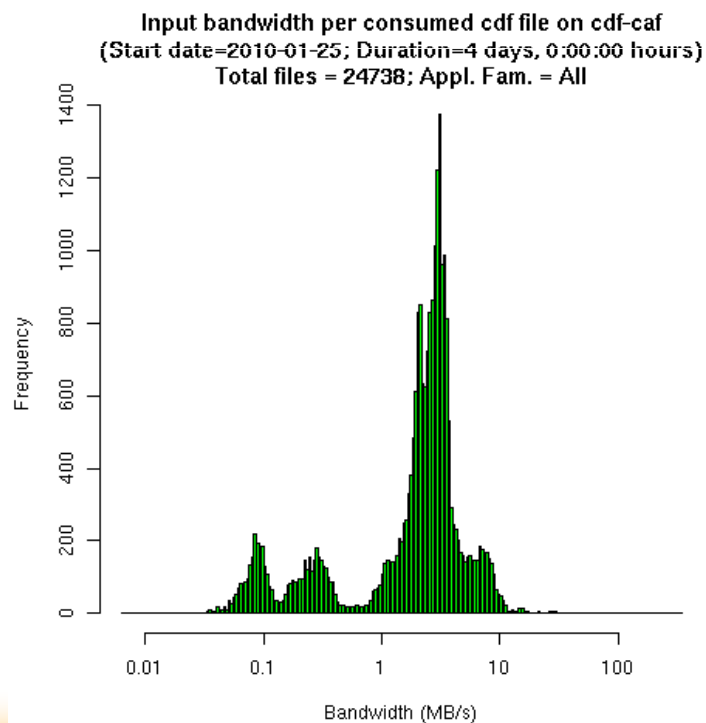
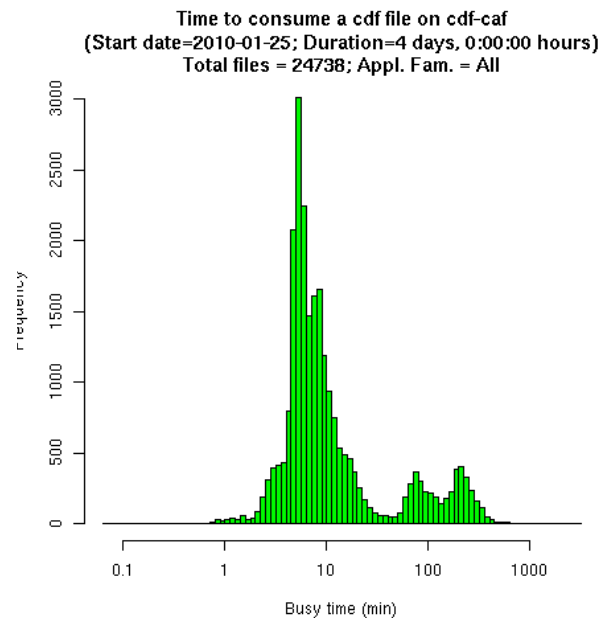
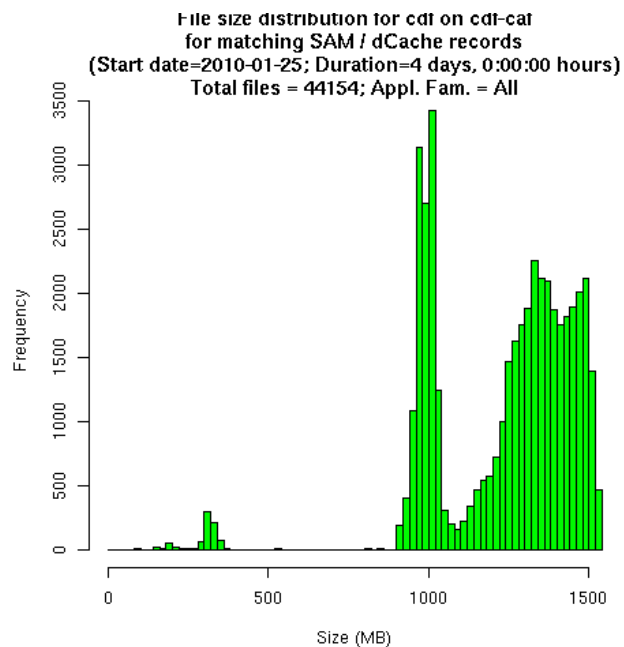
DZero Sam File Access Metrics

- Metrics from the SAM database (for stations FNAL-CABSRV1 and FNAL-CABSRV2)
- Analysis of processes and files from all SAM projects *ended* on the week of Jan 25, 2010
- Jobs run on the CAB cluster. They request file delivery to the SAM-Grid system through the two stations above. Upon request, each system delivers files to a disk cache local to the node where the job is running. Some files may be already in the local cache.
- Analysis code evolved from SAM Monitoring (by Robert Illingworth)
- For all jobs on the CAB cluster, the data shows these distributions
 - file size
 - time to "consume" the file
 - bandwidth from local disk required to read the file.
 - file delivery time
 - number of total and idle processes.
- Details at <http://home.fnal.gov/~garzogli/storage/dzero-sam-file-access.html>



CDF Sam / dCache File Access Metrics

- Metrics from the SAM database and dCache billing log.
- Data from Jan 25, 2010 to Jan 28, 2010.
- The analysis associates event records according to this processing model:
 - A user analysis (client) requests a file to SAM
 - SAM returns a file location in dCache (file *create time* record is created in the SAM DB)
 - The client makes a transfer request to the dCache system (*request* record is logged in the dCache billing log)
 - When the file is available on disk, the client can copy it locally or stream it from dCache (at the end of this process, a *transfer* record is logged in the billing log)
 - The user analysis processes the file and requests a new one to SAM (the status of the file just processed is changed to consumed in the SAM DB and the *update time* for the file is created).
- This analysis is limited in scope by design
- Details at <http://home.fnal.gov/~garzogli/storage/cdf-sam-file-access-per-app-family.html>



Future work

- Evolve test bed
 - Investigate network bottlenecks between client and server
 - Install and run clients on FC; Configure MDS as RAID 10
 - Improve machine monitoring (e.g. install ganglia, HP Lustre monitoring, ...)
 - Test luster with NO stripping (easier recovery)
- Standard benchmarks
 - Run with different configuration e.g. more clients, O(1M) files, different number of OSS, ...
- App-based benchmarks
 - Include more applications: Nova (reco, daq2root, pre-scaled skim), Dzero / CDF (via SAM station)
 - Evolve measurements (study wall vs. real vs. sys time; include CPU, mem, net usage, ...)
 - Collaborate with HEPIX Storage Group
- Move to Phase II and III AND test fault tolerance
- Study Hadoop

Conclusions

- We are about 3 months late according to the original plan: we were too optimistic in the delivery date of FC. Est. completion date: 1st Q 2011.
- We have deployed a Lustre test bed and we are taking measurements
- Thank you for the tremendous response to our requests for help

EXTRA SLIDES

Storage evaluation metrics

Metrics from Stu, Gabriele, and DMS (Lustre evaluation)

- Cost
- Data volume
- Data volatility (permanent, semi-permanent, temporary)
- Access modes (local, remote)
- Access patterns (random, sequential, batch, interactive, short, long, CPU intensive, I/O intensive)
- **Number of simultaneous client processes**
- **Acceptable latencies requirements (e.g for batch vs. interactive)**
- **Required per-process I/O rates**
- **Required aggregate I/O rates**
- **File size requirements**
- Reliability / redundancy / data integrity
- Need for tape storage, either hierarchical or backup
- Authentication (e.g. Kerberos, X509, UID/GID, AFS_token) / Authorization (e.g. Unix perm., ACLs)
- User & group quotas / allocation / auditing
- Namespace performance ("file system as catalog")
- Supported platforms and systems
- Usability: maintenance, troubleshooting, problem isolation
- Data storage functionality and scalability